

## BIROn - Birkbeck Institutional Research Online

Shiode, Shino and Shiode, N. (2020) A network-based scan statistic for detecting the exact location and extent of hotspots along urban streets. *Computers, Environment and Urban Systems* 83 (101500), ISSN 0198-9715.

Downloaded from: <https://eprints.bbk.ac.uk/id/eprint/31871/>

*Usage Guidelines:*

Please refer to usage guidelines at <https://eprints.bbk.ac.uk/policies.html>  
contact [lib-eprints@bbk.ac.uk](mailto:lib-eprints@bbk.ac.uk).

or alternatively

# **A Network-based Scan Statistic for detecting the exact location and extent of hotspots along urban streets**

**Abstract.** Socio-economic activities and incidents such as crimes and traffic accidents have a negative impact on our society, and their reduction has been a priority in our social-science endeavour. These events are not uniform in their occurrences but, rather, manifest a distinct set of concentrations, commonly known as hotspots. Detecting the exact extent, shape and changes in these hotspots can lead to deeper understanding of their cause and help reduce the volume of incidents, yet accuracy of the analytical outcomes using existing methods are often hampered by their reliance on Euclidean distance. This paper proposes a new type of cluster detection method for identifying significant concentration of urban and social-science activities recorded at the individual street-address level. It extends Scan Statistic—a regular hotspot detection method originally developed in the field of epidemiology—by introducing flexible search windows that adapt to and sweep across a street network. Using a set of synthetic data of crime incidents as an example, performance of the proposed method is measured against that of its conventional counterparts. Results from the performance tests confirm that the proposed method is more accurate in detecting the exact locations of hotspots without over- or under-representing them, thus offering an effective means to identify problem places at the individual street-address level. The simulation also demonstrates how well the proposed method captures changes in the intensity of hotspots, which is also something existing methods have struggled with. An empirical analysis is carried out with data on drug, burglary, robbery, as well as thefts from vehicles in Chicago. The study demonstrates the capacity of the proposed method to extract the detailed profile of the concentration of each crime type, which offers interesting insights into their micro-scale patterns which were previously not available at such a fine spatial granularity.

**Keywords:** cluster detection, crime analysis, micro-scale analysis, spatial statistics, street networks

## **1. Introduction**

Cluster detection is a statistical approach used for extracting concentrations of events. It has seen a wide range of applications in recent years, ranging from the agglomeration of suburban communities (Helbich, 2011) and clustering of obesity (Dahly, 2011) to spatial analysis of corner kick goals (Schmicker, 2013), shark attacks (Amin et al., 2012), crime incidents (Nakaya & Yano, 2010; Grubestic et al., 2014) and forest fires (Tuia et al., 2008). In many cases, they utilise the theory and methodology originally developed in the field of epidemiology for identifying concentrations of elevated risks, outbreaks of infectious disease, and for monitoring emerging risks. In particular, Kulldorff's Scan Statistic (Kulldorff & Nagarwalla, 1995; Kulldorff, 1997; Kulldorff, 2009) and its variants, Rushton's spatial filtering (Ozdenerol et al., 2005) and Bayesian disease mapping (Aamodt et al., 2006), have been widely used for analysing the concentration of diseases and health concerns (Osei & Duker, 2008;

Luquero et al., 2011; Desjardins et al., 2018; Li et al., 2019) as well as more prevalent diseases such as various types of cancer (Jemal et al., 2002; DeChello & Sheehan, 2007; Henry et al., 2009).

Scan Statistic (Kulldorff et al., 2003; Song & Kulldorff, 2003) belongs to a group of methods that identifies statistically significant clusters using a search window. It is considered to have overcome two main challenges shared by many other methods, and these are:

- (1) Use of a small number of search window size to reduce the computational load at the cost of reduced accuracy, which tend to induce either undercutting or overshooting the actual clusters. Search windows used in Scan Statistic are flexible in size, thus offering generally more accurate outcomes.
- (2) multiple testing problems that increases Type I errors. For instance, GAM (Geographical Analysis Machine) (Openshaw et al., 1987) and Besag and Newell's method (Besag & Newell, 1991) are known as the two pioneering efforts in this field, both of which use a discrete search window—fixed circle radii in the case of GAM, and a fixed number of cases within a circle adopted by Besag and Newell—and neither adjusts for multiple testing. In contrast, methods such as CEPP (Cluster Evaluation Permutation Procedure) (Turnbull et al., 1990) and Scan Statistic (Kulldorff & Nagarwalla, 1995) use a single unified significance test and are free of multiple testing problems. While CEPP adopts discrete values for population size in each search window, Scan Statistic uses a search window that can continuously change its size.

For these reasons, Scan Statistic is often considered as the most robust and effective cluster detection method. Yet there are still instances where Scan Statistic is found less effective. In the field of urban geography, many applications (e.g. retail and housing hotspots, concentrations of bike traffic accidents, crime events) are confined to street networks (Okabe & Sugihara, 2012); i.e. distribution of events that occur in urban space is affected by the shape and the structure of the street network, and this tendency becomes more pronounced for micro-scale analyses (Braga et al., 2010; Weisburd et al., 2012; di Bella et al., 2017). In such cases, conventional planar methods such as Scan Statistic fail to successfully identify the exact shape of clusters along street networks (Kulldorff, 1997). Recent studies have tried to incorporate certain aspects of street networks into their analysis (Okabe & Okunuki, 2001; Flahaut et al., 2003; Spooner et al., 2004; Yamada & Thill, 2004; Ramp et al., 2005; Okabe et al., 2006; Langen et al., 2007; Maheu-Giroux & de Blois, 2007; Shiode, 2008; Shiode, 2011; Shiode & Shiode, 2009; Shiode & Shiode 2013; Shiode et al., 2014). Among these studies, Shiode (2011), Shiode and Shiode (2013), and Shiode et al. (2014) propose a series of search-window-type methods with the benefit of incorporating the structure of street networks into their analyses. Their studies define a

search window as a flexible sub-network that sweeps along a street network to capture clusters formed along the streets. Shiode (2011) analyse the spatial concentration of events on streets, and Shiode and Shiode (2013) and Shiode et al. (2014) extend it to detect the pattern of space-time concentration of events along street networks. Through simulation and empirical analyses, their studies demonstrate that network-based cluster detection methods offer more effective and accurate alternatives to their conventional counterparts that assume the study area to be a planar, Euclidean space.

The methodological framework proposed by Shiode et al. (2014) is different from that of Scan Statistic in that (1) they use search windows that are flexible in shape but discrete in their size increment, and (2) they apply a single composite hypothesis test for each search window to detect clusters. Use of a discrete window size helps reduce the high computational load for carrying out cluster detection in network space, but it can lead to over-representation of clusters. Similarly, multiple testing of hypotheses could increase Type I error. While recent studies such as the false-discovery-rate correction (Benjamini & Hochberg, 1995) propose to correct for multiple testing, they are subject to ongoing debates. Instead of bringing them into its framework, this paper develops a network equivalent of Scan Statistic for urban applications by utilising a flexible search window. Developing a variant of Scan Statistic for detecting flexible shaped clusters is not a new idea (Duczmal & Assunção, 2004; Patil & Taillie, 2004; Tango & Takahashi, 2005; Kulldorff et al., 2006), and some have indeed outperformed the standard circular Scan Statistic. However, none of them are designed specifically for analysing events recorded along street networks of an urban environment.

## 2. Methodology

### *Network-based Scan Statistic (NetScan)*

Against this background, this paper proposes a network version of Scan Statistic, *NetScan* and examines its validity in the context of urban crimes. The original Scan Statistic (Kulldorff & Nagarwalla, 1995) is designed to detect regions in which the number of observed incidents is significantly higher than expected. For each instance of searches, the likelihood ratio is computed by counting the number of observed incidents inside and outside that window. The window that maximises the likelihood ratio statistic becomes the most likely cluster. The shape of a search window can take circular, elliptic and other irregular shapes (Duczmal et al., 2006) but they have, so far, always used a search window with an area. This paper departs from that approach by replacing a standard planar search window with a network-based search window, or a set of connected line segments that

extends from a reference point on a street network with its total length ranging from zero to a pre-determined maximum length; and uses it to sweep across the entire extent of the street network in the study area.

The original Scan Statistic usually assumes inhomogeneous Poisson process with an intensity proportional to a known attribute such as population (Kulldorff, 1997). It aggregates the incident counts to a finite number of areal units such as the administrative district areas, typically represented by their centroids. Each location is assigned the respective associated cases (e.g. the number of observed incidents within the area) and a base population (e.g. *at-risk* population). The main interest lies in detecting clusters that cannot be explained by the baseline process (Kulldorff, 1997).

While most studies using Scan Statistic adopt *inhomogeneous Poisson process* as their underlying process, this paper assumes *continuous and homogeneous Poisson process* where observations are distributed randomly and continuously throughout the street network of a study area with a constant intensity according to a homogeneous Poisson process. It serves the scope of this paper well, as estimating the base population for every location along a street network would be impractical, and the notion of at-risk population is less essential for interpreting the types of crime covered in this paper. Using this assumption, observations can be recorded at any location along the street network. They are then searched with a flexible search window for any clusters of crime incidents that are unlikely to emerge, if each incident occurred independently and randomly (Kulldorff, 2009).

#### *Construction of network-based search windows*

The process of detecting a statistically significant cluster with NetScan is as follows. Let  $N$  be the street network in a study area (hereafter called study network  $N$ ), and  $n_G$  be the total number of observed points found on  $N$ . Let  $\mathbb{Z}$  be the collection of network-based search windows created on  $N$ ,  $Z$  be an element of the set  $\mathbb{Z}$ , and  $n_z$  be the observed number of points inside  $Z$ . Under the null hypothesis, no cluster is assumed to exist on  $N$ , and the number of incidents in each search window is Poisson distributed with expected values proportional to the length of that search window. In other words, the expected number of incidents for window  $Z$  can be defined as  $\lambda(Z) = n_G |Z| / |N|$ , where  $|Z|$  is the length of the search window and  $|N|$  is the length of the entire study network. Conversely, the alternative hypothesis  $H_1(Z)$  assumes the presence of a hotspot in region  $Z$  and that the expected counts inside and outside of a window  $Z$  are multiplied by unknown constants  $p$  and  $q$  respectively, where  $p > q$ . In this case, the maximum likelihood estimate of  $p$  is  $C_{in} / A_{in}$  and the maximum likelihood estimate of  $q$  is

$C_{out}/\Lambda_{out}$ , where  $C_{all}$  and  $\Lambda_{all}$  are the aggregate observed count  $\sum c_i$  and expected count  $\sum \lambda_i$  for all line segments  $z_i$  ( $z_i \in Z$ ), and  $C_{out} = \sum_{z_i \notin Z} c_i$  and  $\Lambda_{out} = \sum_{z_i \notin Z} \lambda_i$ , respectively.

$H_0$ :  $c_i \sim \text{Poisson}(\lambda_i)$  for all line segments  $z_i$ .

$H_1(Z)$ :  $c_i \sim \text{Poisson}(p\lambda_i)$  for all lines segments  $z_i$  in  $Z$ , and  $c_i \sim \text{Poisson}(q\lambda_i)$  for all line segments  $z_i$  outside  $Z$ , for some constant  $p > q$ .

Our goal is to determine whether large observation counts within a window is due to chance fluctuations, and to test if this specific  $Z$  constitutes a statistically significant cluster. Based on these hypotheses, NetScan detects increased count in a region  $Z$  if the ratio of the observed to the expected counts (calculated as a form of a likelihood ratio) is higher inside than outside the region. NetScan uses a similar approach to the standard Scan Statistic for deriving the likelihood ratio for the continuous homogeneous Poisson model on a street network.

Suppose that  $L(Z, p, q)$  is the likelihood function, and the likelihood ratio  $T$  is the likelihood under the alternative hypothesis  $L(Z)$  divided by the likelihood under the null hypothesis  $L_0$ . It shows how likely the observed data for  $Z$  are given a differential rate of incidents inside and outside of  $Z$ . The likelihood ratio  $T$  can be written as

$$T = \frac{L(Z)}{L_0} = \frac{\sup_{Z \in \mathbb{Z}, p > q} L(Z, p, q)}{\sup_{p=q} L(Z, p, q)} = \sup_{Z \in \mathbb{Z}} \left( \frac{n_Z}{\lambda(Z)} \right)^{n_Z} \left( \frac{n_G - n_Z}{\lambda(G) - \lambda(Z)} \right)^{n_G - n_Z}$$

if there is at least one zone such that  $\frac{n_Z}{\lambda(Z)} > \frac{n_G - n_Z}{\lambda(G) - \lambda(Z)}$ , (1)

and  $\lambda=1$ , otherwise.

By maximising the likelihood ratio over all  $Z$ , the single  $Z$  that forms the most likely cluster can be identified. The likelihood ratio calculated for this window constitutes the maximum likelihood ratio test statistic, i.e. the hotspot most unlikely to be found under the null hypothesis. The maximum likelihood ratio can be derived as

$$T = \max T(Z) = \max_{Z \in \mathbb{Z}} \left\{ \frac{L(Z)}{L_0} \right\} \quad (2)$$

where window  $\hat{Z} \in \mathbb{Z}$  is derived as the most likely cluster, and  $\hat{Z}$  is selected to satisfy  $L(\hat{Z}) \geq L(Z)$  for all  $Z \in \mathbb{Z}$ . Once the region with the maximum likelihood is found, it is examined for its statistical significance; i.e. if the  $p$ -value for  $\hat{Z}$  is less than some fixed significance level  $\alpha$ .

Theoretically speaking, the null hypothesis distribution of the maximum likelihood ratio test statistic is intractable, but it can be approximated with Monte Carlo simulations. It gives us the  $p$ -value for each cluster, from which we can evaluate its statistical significance (Duczmal et al., 2006). The method is known as randomisation testing and uses a large number  $B$  (usually 999) of random replications of the data set under the null hypothesis. For each replication, the maximum likelihood ratio  $T^*$  is calculated (Duczmal & Assunção, 2004). The statistical significance of the most likely hotspot can be calculated by comparing  $T(Z)$  to these replica values of  $T^*$ . The  $p$ -value of region  $Z$  can be computed as  $\frac{B_{beat} + 1}{B + 1}$ , where  $B$  is the total number of replicas produced, and  $B_{beat}$  is the number of replicas with  $T^*$  greater than  $T(Z)$ . If this  $p$ -value is less than the significance level  $\alpha$ , we conclude that the region forms a significant cluster.

Results from Monte Carlo simulations can be also used for detecting other clusters. By comparing the  $p$ -value of each potential cluster, all hotspots with a  $p$ -value less than the significance level  $\alpha$  are reported as secondary clusters. The computational effort for PL-Scan and NetScan are comparable in that both are executed against a similar number of seed points with the same Monte-Carlo runs. The difference in the computational time arises from the shortest-path search for constructing network search windows which, in a typical urban environment, would come to  $O(m + n \times \log(n))$ , where the number of edges ( $m$ ) and that of nodes ( $n$ ) are finite. It is also worth noting that the very purpose of conducting micro-scale analysis is to understand hotspots at a local scale, which would not exceed the bounds of a single city.

### 3. Detection of micro-scale hotspots using synthetic point data

#### *Synthetic clustered point patterns: the Poisson cluster process*

Using a set of synthetic points generated along a small street network with known clusters, this section carries out a simulation study and compare the performance of NetScan and its conventional counterpart, the standard Spatial Scan Statistic (hereafter PL-Scan) in accurately detecting hotspots. A proprietary programme code for NetScan is prepared to facilitate the simulation study. The conventional methods of circular and elliptic PL-Scan are carried out using the SaTScan software.

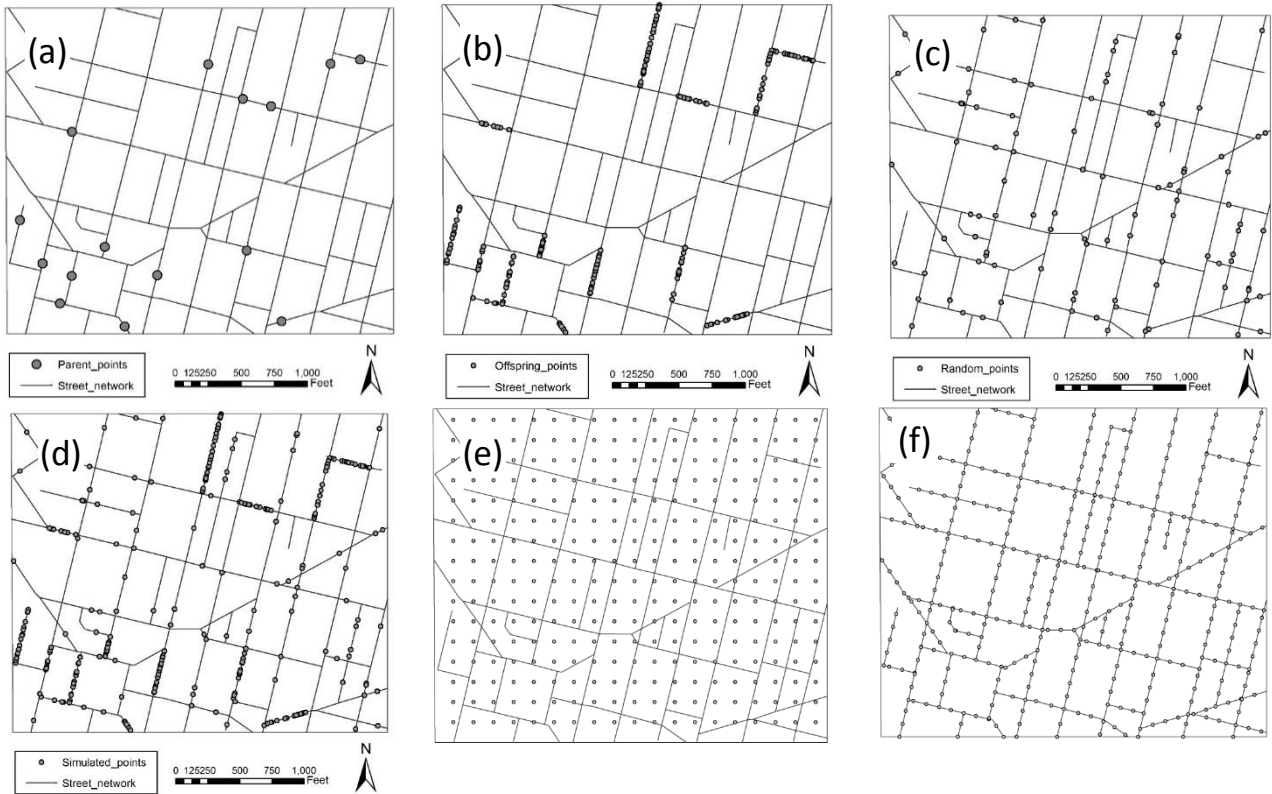
A number of point processes are used in the literature to generate the characteristic cluster patterns, including Cox Process, Neyman type A process, Negative binomial process and Poisson cluster process (Cliff & Ord, 1981; Boots & Getis, 1988; Cressie, 1991). This paper adopts *the Poisson cluster process*

(Upton & Fingleton, 1985; Diggle, 2003) and extends it to the network space to create clusters of points at locations of choice. The process consists of three steps:

Step 1: *Generate a set of points called parent points (Diggle, 2003), around which the clusters will be injected.* This is realised by randomly placing a predetermined number of points on  $N$ .

Step 2: *Generate a set of points called offspring points (Diggle, 2003) around each parent point.* In other words, we add a set of points on a line segment that contains at least one parent point to create a cluster around each parent point. The line segments are bounded respectively by their two end points, and a fixed number of offspring points are randomly placed between them.

Step 3: *Generate a uniform random distribution on  $N$  and superimpose them on the offspring points from Step 2.* This produces a hypothetical distribution of street crime incident locations.

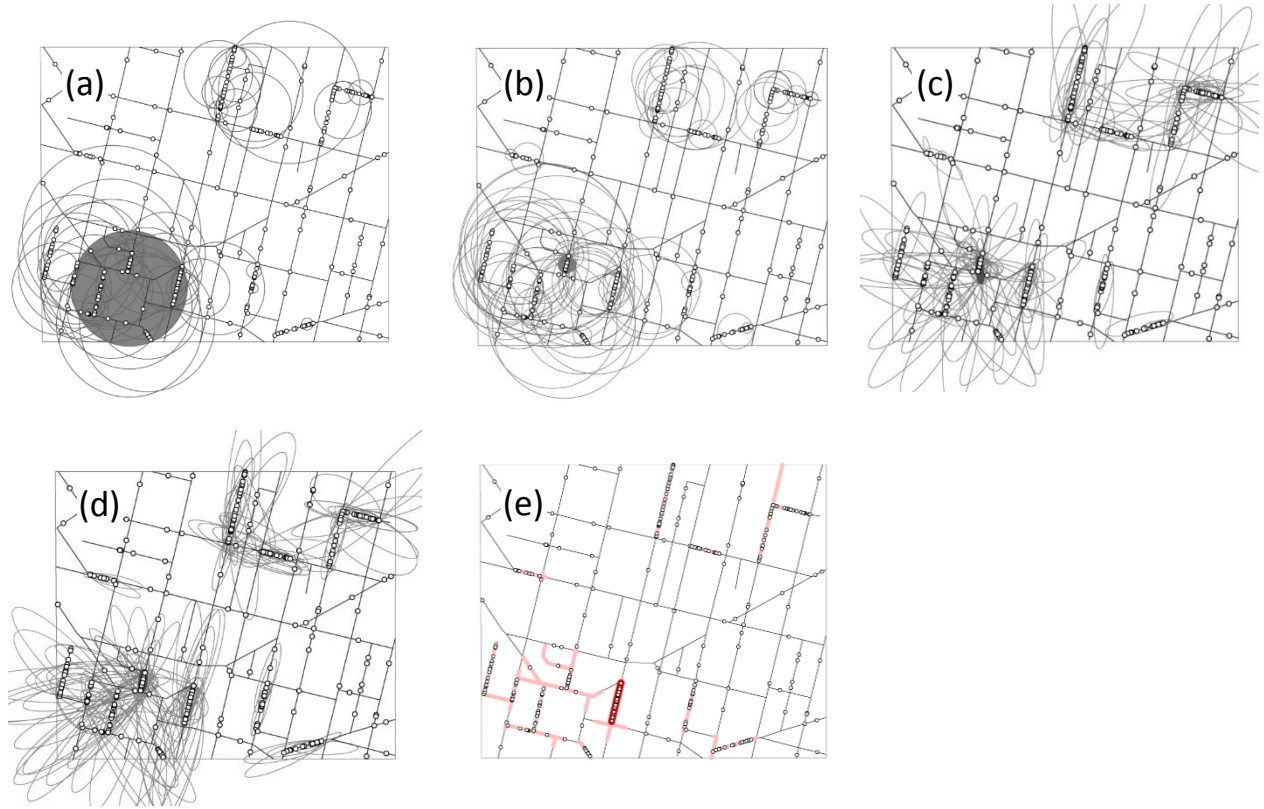


**Figure 1.** Dataset for simulation testing created along the streets of downtown Buffalo, NY: (a) 15 random parent points on  $N$ ; (b) 200 random offspring points on the 14 line segments identified by at least one parent point; (c) 100 random points on  $N$ ; (d) a synthetic Poisson clustered distribution of 300 points consisting of 200 random offspring points and 100 random points on  $N$ ; (e) 304 grid-type reference points placed across the study area, and (f) 394 network reference points placed along the street network.



The simulation data covers an area of 900m by 750m from downtown Buffalo, NY with the streets extending to over 12,500m (Figure 1). In this example, 15 parent points are randomly added for selecting the line segments with clusters (Figure 1(a)). Two are on the same segment, resulting in 14 segments for generating points. A total of 200 offspring points are randomly placed across these segments with a consistent density (Figure 1(b)). Finally, 100 random points are added across  $N$  (Figure 1(c)) to complete an inhomogeneous clustered pattern of 300 points (Figure 1(d)).

To explore the sensitivity of PL-Scans to the configuration of street networks, two sets of reference points are prepared: (1) grid-type reference points spread across the study area (Figure 1(e)), and (2) reference points placed along the street network (Figure 1(f)). The grid-type reference points are placed at an interval of 150ft (45.72m), creating a 19-by-16 grid of 304 points; while the network-based reference points (Figure 1(f)) are aligned with the street network, comprising 394 points at an interval of approximately 100ft (30.48m).



**Figure 2.** The most likely cluster (shown in dark-grey shade) and other statistically significant clusters (shown as empty circles or ellipse) among the 300 observed points detected by (a) circular PL-Scan using the grid-type reference points; (b) circular PL-Scan using the network-based reference points; (c) elliptic PL-Scan using the grid-type reference points; (d) elliptic PL-Scan using the network-based reference points; and (e) the most likely cluster (shown in red) and other statistically significant clusters (shown in pink) detected by NetScan.

### *Applying PL-Scan to the synthetic data*

Figures 2(a) and 2(b) show the results from the application of circular PL-Scan to the synthetic point data using the two sets of reference points. The circular areas in dark-grey shade mark the most likely cluster under each search. Other empty circles represent all secondary clusters that were also detected as significant. While the distinction between the most likely cluster and secondary clusters may have a substantive meaning in epidemiological contexts, it becomes less crucial when detecting problem places in an urban context such as crime hotspots, as all cluster locations require as much attention as other cluster locations do. In this paper, the most likely cluster is being highlighted as an illustrative example of those detected by each method.

Figure 2(a) shows that the application of circular PL-Scan using grid-type reference points could result in the most likely cluster with a relatively wide area that spans across multiple cluster locations (with a radius of 486 ft,  $p$ -value=0.001 to cover 83 points across 5 parent points), whilst also detecting 46 secondary clusters as significant.

In contrast, application of circular PL-Scan with reference points placed along the street network results in the most likely hotspot being confined to a smaller circle consisting of only 13 points cluster locations around a single parent point (with a radius of 46.73ft,  $p$ -value=0.001) whilst 72 other clusters are also detected as statistically significant. (Figure 2(b)).

The difference in the patterns of detected clusters between the two sets of reference points is mainly in the number of clusters detected, and their overall distribution is quite similar. Use of reference points along street networks does allow a slightly more focused identification of a cluster in that some of the larger clusters detected with grid-type reference points are eliminated or reduced in size. However, many of the clusters still tend to over-represent the actual cluster by merging two or more clusters into one, or covering points outside those clusters as an integral part.

#### *Applying elliptic PL-Scan to the synthetic data*

In order to examine if the accuracy of cluster detection along street networks can be improved by changing the shape of search windows, an elliptic variation of Scan Statistic (hereafter referred to as elliptic PL-Scan) is also applied to the same dataset of synthetic points. Elliptic PL-Scan is similar to circular PL-Scan, except it uses an elliptic search window (Kulldorff, 2006).

Figures 2(c) and 2(d) show the results from the application of elliptic PL-Scan to the synthetic point data using the two sets of reference points. The ellipse in dark-grey shade marks the location detected as the most likely cluster under each search. Both cases identify the same hotspot location that

covers a small cluster formed around the shortest line segment of all cluster locations. This tendency was also confirmed by circular PL-Scan when using reference points along the street network (Figure 2(b)).

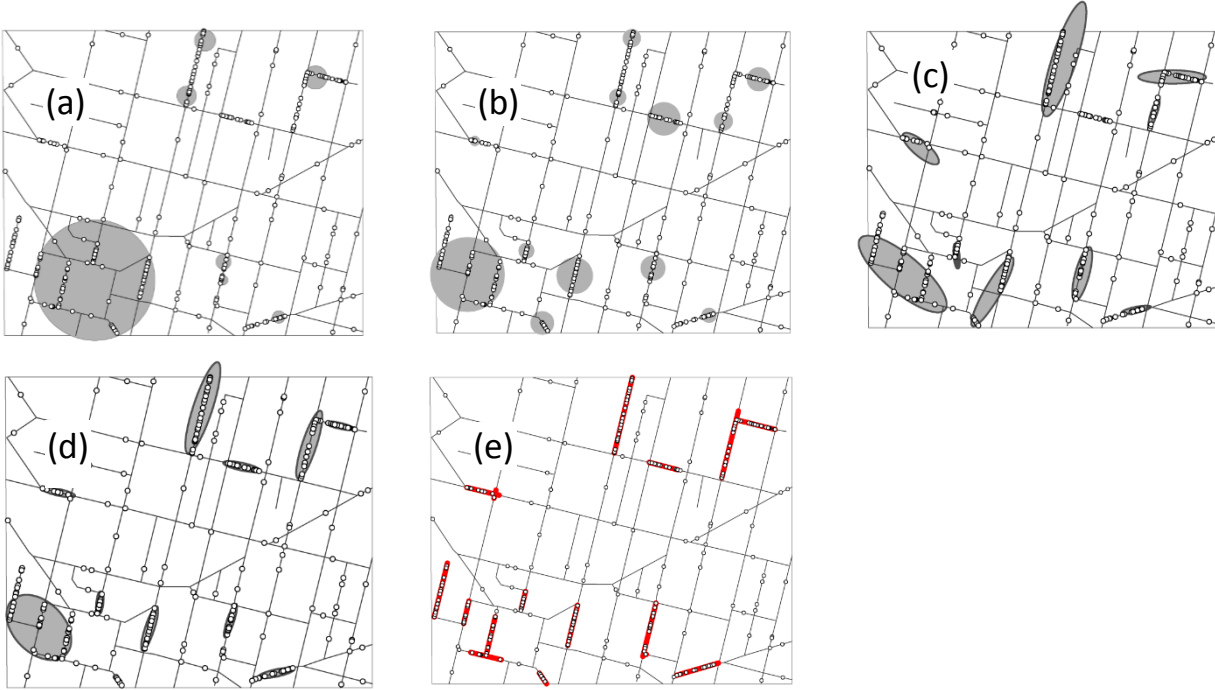
While the use of reference points along the street network results in detecting a greater number of overlapping clusters, their overall patterns are similar, showing some tendency to over represent cluster locations by either covering excess area or incorrectly merging two or more cluster locations into a single cluster.

#### *Applying NetScan to the synthetic data*

Finally, NetScan was applied to the same set of synthetic point distribution with 300ft (91.44m) as the maximum threshold for the size of network-based search windows. Figure 2(e) shows the most likely cluster (shown in red) and other overlapping clusters (shown in pink) detected along the street network by NetScan. Each of these clusters covers exactly one cluster location identified by one of the 14 line segments containing at least one parent point. The line segment detected as the most likely hotspot by NetScan encloses 22 points ( $p\text{-value}=0.001$ ).

Comparing Figures 2(a)-2(e) brings out some interesting difference between the clusters detected by circular PL-Scan and elliptic PL-Scan (using either grid-type reference points or those along street networks), as well as NetScan using continuous search. The most notable difference is that both circular and elliptic PL-Scans tend to over-represent cluster locations, whereas NetScan pinpoints the cluster locations without notable exaggeration of the extent of clusters. In terms of the most likely cluster, NetScan detects a different hotspot to that detected by the other two. This is because, in theory, the density of points should be consistent across all 14 line segments; however, in reality, these randomly assigned points will have created some degree of unevenness of point density within and between each line. It affects calculation of the likelihood ratio, especially as the existing range of methods calculate the likelihood ratio for a circular or an elliptic area, as opposed to pinpointing locations along the street network.

The difference between the hotspots detected by PL-Scans and those by NetScan owes mainly to the shape of the PL-Scan search windows. To capture the entire extent of a cluster detected by NetScan, circular PL-Scan needs to cover the convex hull of the cluster, which increases the window size as well as the expected number of counts, thus reducing the likelihood ratio.



**Figure 3.** Non-overlapping clusters among the 300 observed points detected by (a) circular PL-Scan using the grid-type reference points; (b) circular PL-Scan using the network-based reference points; (c) elliptic PL-Scan using the grid-type reference points; (d) elliptic PL-Scan using the network-based reference points; and (e) NetScan.

Locations of other significant clusters detected by each method are broadly similar across all methods. However, many of the clusters detected with PL-Scans combine two or more cluster locations, whereas NetScan pinpoints individual clusters without merging them with others. Also, clusters detected by NetScan are, by definition, confined to the street network, which minimises any overspill or over-representation of a cluster. This becomes evident when comparing the total area of clusters detected by each method. Circular and elliptic PL-Scans cover roughly 30% and 20% of the study area, respectively, regardless of the patterns of the reference points used.

### *Detecting non-overlapping clusters*

To improve the clarity of the cluster representation, non-overlapping clusters, or the cluster with the highest likelihood ratio amongst overlapping clusters, were also extracted. Figures 3(a) and 3(b) show the results from the circular PL-Scan searches for non-overlapping clusters. Outcomes derived with the two sets of reference points are broadly similar, but the reference points assigned along the street network yields a better performance, successfully detecting all cluster locations with less over-representation (Figure 3(b)). Even then, clusters detected with circular PL-Scan tend to over-represent the larger clusters (i.e. a large concentration of points tend to be shown with a cluster that covers a

wider area than it should) and under-represent smaller concentrations (i.e. detecting only a portion of a small cluster).

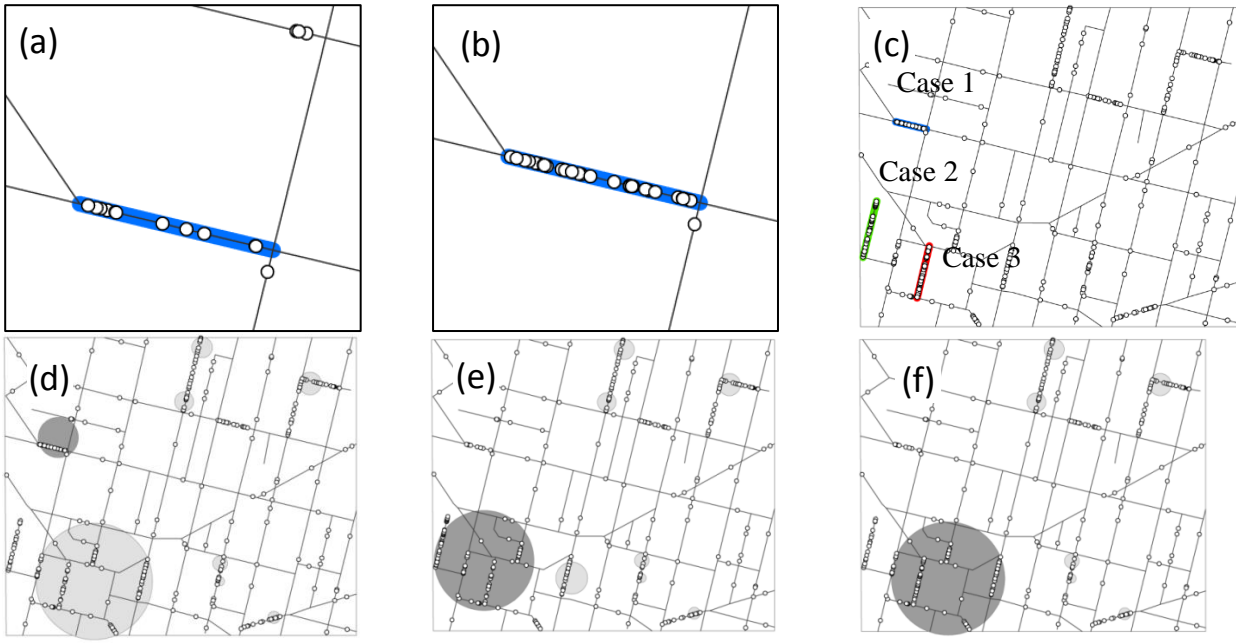
Figures 3(c) and 3(d) show the results from the elliptic PL-Scan searches for non-overlapping clusters using the same set of synthetic points. As is the case with circular PL-Scan searches, the grid-type reference points tend to produce clusters that span across more than one hotspot, whilst some other cluster locations remain undetected. In contrast, reference points assigned along the street network improves the performance of elliptic PL-Scan, allowing it to successfully detect all cluster locations. It confirms that, for detecting clusters along a street network, reference points that are also assigned along street networks would yield better performance, regardless of the shape of search windows used. In fact, aside from one case of over-representation of a cluster, using elliptic PL-Scan with network-based reference points (Figure 3(d)) gives quite an accurate impression of the cluster locations.

Figure 3(e) shows the non-overlapping clusters detected through a spatially continuous search with NetScan using a flexible search window. Eliminating the overlapping clusters helps illustrate more clearly that NetScan accurately detects all cluster locations around the respective parent points. While PL-Scans identify most of these locations, they tend to over- or under-represent the extent of clusters (Figures 3(b) and 3(d)). NetScan not only identifies the cluster locations correctly but also detects the exact extent of spatially significant clusters with very little over-representation.

#### **4. Sensitivity of PL-Scans and NetScan**

Comparing the performance of PL-Scans and NetScan with a fixed set of points only shows their performance under static state, and not how they respond to changes in the locations of high-risk areas. In order to examine their sensitivity to such changes, 20 new points are added to one of the secondary clusters, to turn it into the most significant cluster.

Figures 4(a) and 4(b) respectively show the state of a cluster before and after new points are added to that location. Figure 4(c) shows three cluster locations used as illustrative examples for adding new points; i.e. in each case, new points were added to one of the three locations, making it the most significant cluster with the highest likelihood ratio. In all three instances, NetScan detected the changes in the density of points at the respective cluster location and correctly identified the new location as the most likely cluster.



**Figure 4.** (a)-(c): Preparing data to explore sensitivity of PL-Scan and NetScan to changes in the most likely hotspot: (a) a cluster location in its original state, (b) the same cluster location with additional 20 points, and (c) illustration of three different clusters, each with additional 20 points; (d)-(f): The updated most likely hotspot detected by PL-Scan (with reference points assigned along the street network): (d) Case 1, (e) Case 2, and (f) Case 3.

Figures 4(d)~4(f) show the results from the application of circular PL-Scan to the same three cases as in Figure 4(c) where new points are injected at the respective location. For Case 1, both circular PL-Scan and NetScan successfully identified the new hotspot as that with maximum likelihood, although the hotspot detected by PL-Scan is slightly off-centred and extends to another line segment with no points (Figure 4(d)). In Case 2, PL-Scan continued to point to the same location that contains the previous hotspot and failed to detect the new hotspot as the most significant one (Figure 4(e)), whereas NetScan successfully detected the new, intensified concentration. This is because the change was subtle in this case where the likelihood ratio for the new hotspot was only slightly higher than the rest, and the p-values remained the same for all detected hotspots. In Case 3, the location of the cluster with the highest density of points was successfully covered by both PL-Scan and NetScan, as both methods moved the location of the most significant cluster. However, the result from PL-Scan covers a larger area and makes it difficult to learn where in particular the risk of crime was elevated within that hotspot location.

Results from the application of the two Scan Statistic methods, namely PL-Scan and NetScan, to both the original and the modified synthetic distributions suggest that: (1) NetScan has the capacity to pinpoint the hotspot location more precisely and concisely than PL-Scan can, and (2) NetScan is more sensitive to changes in the density of points, and detects even a subtle change.

## 5. Performance Assessment

Discussion of the detection accuracy in the preceding section was based on a single set of synthetic point distribution, which is not sufficient for understanding the variability of possible outcomes. This section examines the detection accuracy by introducing and systematically assessing nine other synthetic distributions and comparing the performance of the three methods for detecting the clusters within.

In the field of syndromic surveillance, outcomes from a surveillance system is often assessed with respect to (1) the *completeness* in detecting all hotspot locations, and (2) the *extent of coverage* of the correct hotspot areas (Nordin et al., 2005). Completeness, or the *detection power*, refers to the success rate in detecting the clusters. Given the scope of this paper, we will measure the accuracy of the coverage of hotspot areas. Specifically, we will follow the approach taken by Forsberg, et al. (2005), Huang, et al. (2007), Takahashi, et al. (2008), and Neil (2009), who combined two intertwined measures in their performance assessment: *positive predictive values (PPV)* and *sensitivity*. The PPV is the proportion of the true “regions” within the detected clusters, while sensitivity measures the probability of detecting a “region” that actually constitutes a cluster (Takahashi et al., 2008). As those studies were carried out in the context of syndromic surveillance, they dealt mainly with data collated to areal units, and the PPV and sensitivity values were based on *the number of regions* detected or covered. As we are focusing on individual places represented by points, this paper examines *the number of reference points* instead.

Let  $S_{true}$  be the *true hotspot regions* (i.e. the actual synthetic hotspots), and let  $S^*$  be the *detected regions* (i.e. the regions detected as statistically significant by the respective method). Using the reference points, the PPV can be defined as the ratio of correctly detected locations (i.e. the number of detected reference points  $r_i$  within the respective  $S_{true}$ ) to all detected locations (i.e. the number of reference points  $r_j$  in the respective  $S^*$ ), which measures the degree of overshooting:

$$PPV = \frac{\#\{r_i \in S^* \cap S_{true}\}}{\#\{r_j \in S^*\}} \quad (3)$$

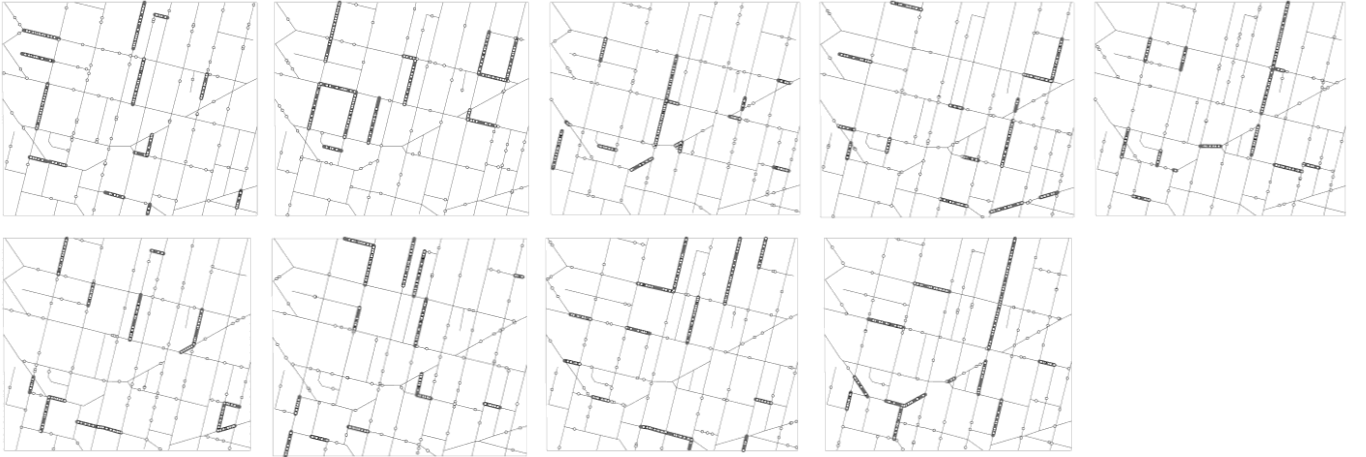
The PPV becomes low when the detection accuracy suffers from a high level of overshooting, and it becomes high when there is less overshooting.

The specificity accounts for the degree of undershooting by comparing the ratio of correctly detected locations (i.e. the number of detected reference points  $r_i$  within the respective  $S_{true}$ ) to all true hotspot locations (i.e. the number of reference points  $r_j$  in the respective  $S_{true}$ ):

$$Specificity = \frac{\#\{r_i \in S^* \cap S_{true}\}}{\#\{r_j \in S_{true}\}} \quad (4)$$

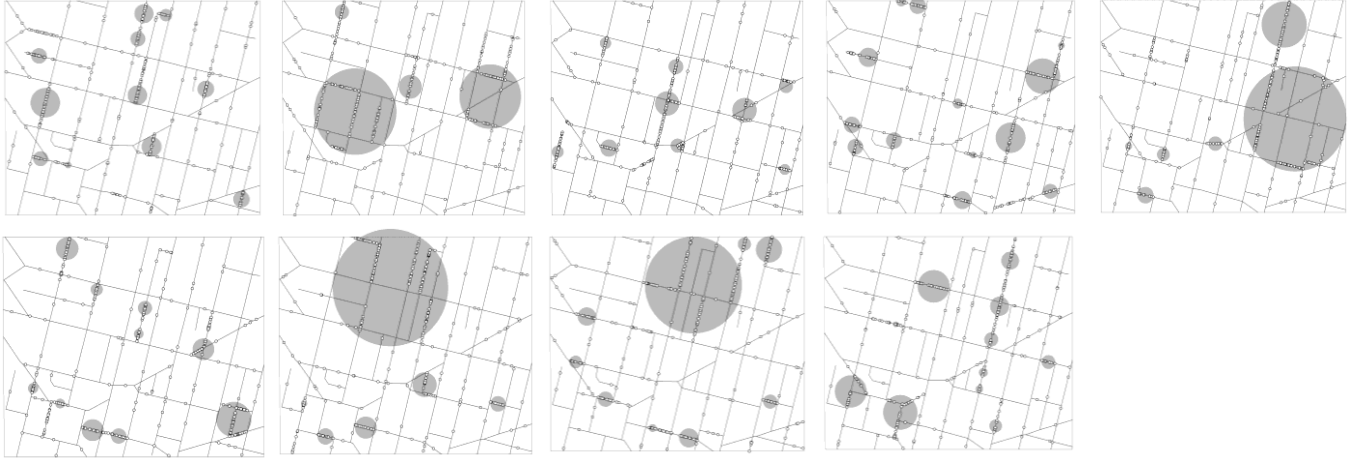
The specificity becomes low when it suffers from a high level of undershooting, and it becomes high when there is less undershooting. For both PPV and specificity, the higher the value, the better the performance of the method examined, with a score of 1.0 indicating no undershooting or overshooting.

In order to test the performance of the proposed method, nine additional realisations of the Poisson cluster model were generated (Figure 5). Figures 6(a), 6(b) and 6(c) show the outcomes of hotspot detection using circular PLScan, elliptic PLScan and NetScan, respectively. All three methods used the same set of reference points assigned along the street network, as applying circular and elliptic PLScans using reference points along the street network returned more accurate outcomes than the results obtained with grid-type reference points.

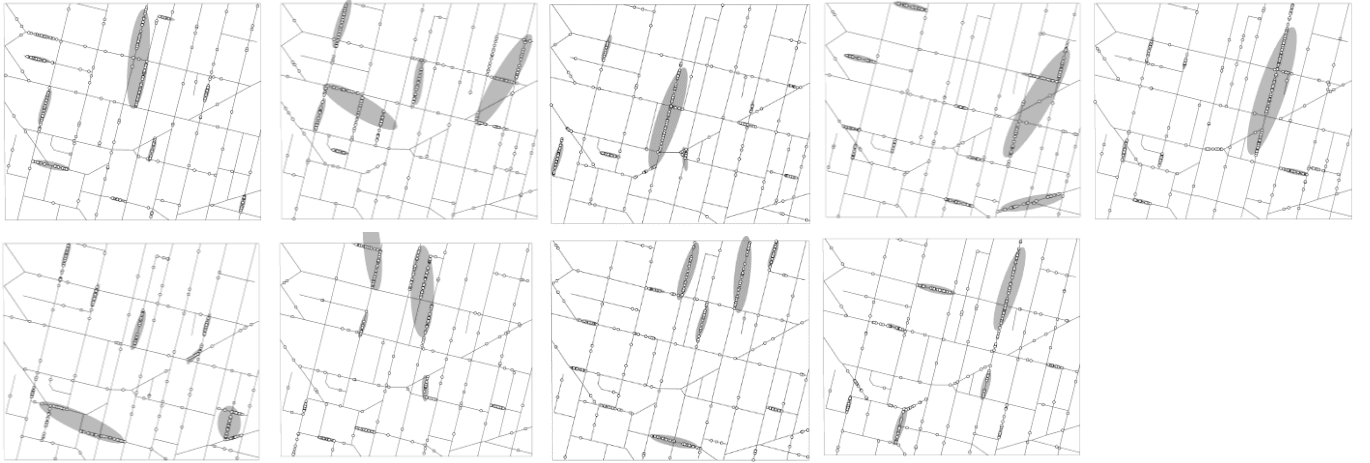


**Figure 5:** Nine realisation patterns of the Poisson cluster process, in which hotspots were injected on 14 line segments with a total number of 300 points (200 cluster points injected on the selected line segments and 100 random points spread across the entire street network).

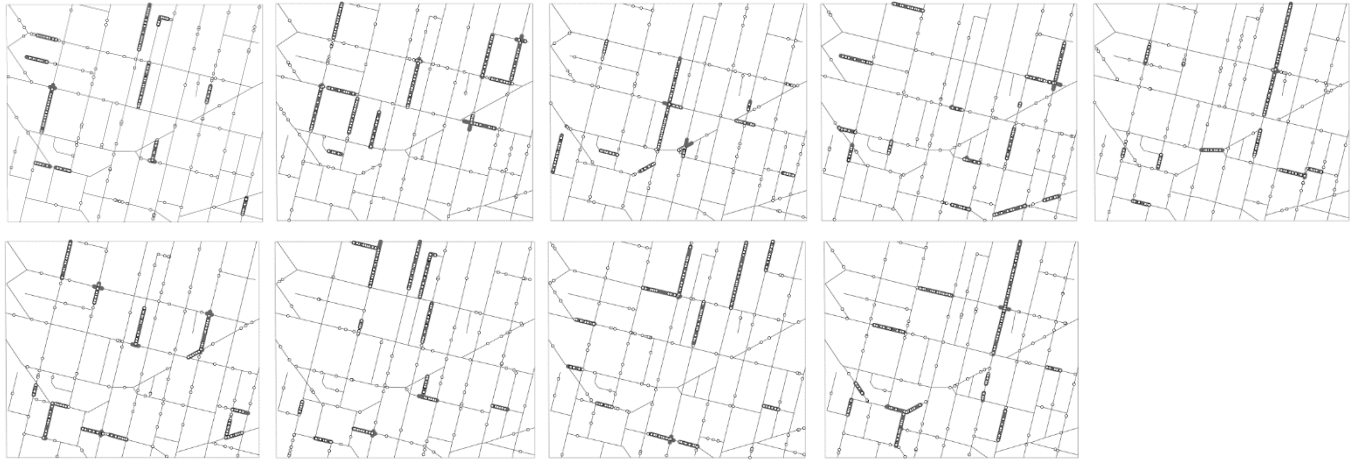




(a)



(b)



(c)

**Figure 6:** Non-overlapping clusters detected for the nine synthetic Poisson cluster patterns from Figure 5, using (a) circular PLScan; (b) elliptic PLScan; and (c) NetScan.

Table 1 shows PPV and specificity calculated for the 10 synthetic Poisson cluster patterns, including the first set from earlier section (Figure 1(d)) and the 9 newly generated patterns (Figure 5). NetScan returns the highest mean values for both PPV and specificity, confirming the visual observation from Figures 3(a), 3(b) and 3(e) in NetScan returns the highest detection accuracy, followed by elliptic PLScan and circular PLScan.

**Table 1:** Performance of circular PLScan, elliptic PLScan and NetScan assessed in terms of their PPV and specificity scores for 10 patterns of synthetic point distributions.

Realised Pattern	Circular		Elliptic		Network	
	PPV	Specificity	PPV	Specificity	PPV	Specificity
1	0.74	0.52	0.93	0.74	1.00	0.93
2	1.00	0.48	0.83	0.76	1.00	0.72
3	0.53	0.60	0.93	0.59	0.96	0.78
4	0.84	0.57	0.84	0.57	0.92	0.78
5	0.92	0.60	0.83	0.77	0.96	0.89
6	0.37	0.56	0.79	0.58	1.00	0.75
7	0.85	0.48	0.82	0.56	0.98	0.86
8	0.54	0.81	0.86	0.69	1.00	0.80
9	0.58	0.64	0.90	0.62	1.00	0.83
10	0.91	0.51	1.00	0.58	1.00	0.79
Mean	<b>0.73</b>	<b>0.58</b>	<b>0.87</b>	<b>0.64</b>	<b>0.98</b>	<b>0.81</b>
Standard Deviation	<b>0.21</b>	<b>0.10</b>	<b>0.06</b>	<b>0.09</b>	<b>0.03</b>	<b>0.07</b>
Coefficient of Variation	<b>0.29</b>	<b>0.17</b>	<b>0.07</b>	<b>0.13</b>	<b>0.03</b>	<b>0.08</b>

The relatively low PPV scores from circular and elliptic PLScans indicates they have overshot, which was due to the shape of their search windows. Circular and elliptic hotspots tend to be inclusive and merge multiple hotspots into a single large “cluster”, which results in a large area with no hotspots (Figure 6(a), 6(b)). The level of overshooting is heavily affected by the spatial arrangement of the events. It increases when more injected hotspots are found near each other, as they are more likely to be mistaken as a single cluster. The level of overshooting becomes particularly high (or the PPV score is low) with circular PLScan when large circular hotspots are formed (e.g. Patterns 3, 6, 8 and 9) to cover a large vacant area outside true hotspots, whereas some other patterns result in a number of small circles with a relatively low level of overshooting (Figure 6(a)). This contrast yields a high variability in the PPV scores, which is also reflected in the high value of coefficient of variation. Elliptic PLScan (Figure 6(b)) puts on a better performance both in terms of the PPV scores and the reasonably low value of coefficient of variation, but is not free of such fluctuation between different patterns. NetScan on the other hand constantly detects each hotspot separately, and this is reflected in its consistently high PPV scores and small variability (Figure 6(c)).

In order to further compare the outcomes from the three methods, Mann-Whitney's U test was carried out (Table 2). It confirmed that the PPV and the specificity scores from NetScan are significantly different from those from circular and elliptic PLSans (at  $\alpha=0.05$ , two-tailed test); whilst the results from circular PLSan and elliptic PLSan were not significantly different.

**Table 2:** Results from the Mann-Whitney's U test for PPV and Specificity.

	Circular PLSan	Elliptic PLSan	NetScan	
Circular PLSan		U=33.0 (p=.210) U= 28.0 (p=.103)	U= 7.5 (p=.002*) U= 6.0 (p=.001*)	PPV
Elliptic PLSan	U=33.0 (p=.210) U= 28.0 (p=.103)		U=9.0 (p=.002*) U= 5.0 (p=.001*)	Specificity
NetScan	U= 7.5 (p=.002*) U= 6.0 (p=.001*)	U=9.0 (p=.002*) U= 5.0 (p=.001*)		

In most cases, PPV scores were generally higher than the specificity scores. High PPV confirms close match with true hotspots and low specificity indicates the true hotspots to be larger than the ones detected which, when combined, suggest that, we generally have more undershooting than overshooting. This could be due to the nature of Monte Carlo simulation. When placing the concentrated points, points are injected randomly across the selected line segments, but there is no guarantee of getting them scattered from one end of the line segment to the other end; i.e. in reality, there may be some vacant or sparse area. In such a case, the extent of a network segment detected by NetScan as a cluster may not necessarily cover the entire line segment. The result is undershooting of hotspots, where points around the edge of the line segment tend to remain undetected. While this does not affect the comparative merit of the relative values of specificity across the different methods, it generally decreases the specificity scores when compared to how all methods perform on their PPV scores.

## 6. Empirical analysis of micro-scale hotspots using NetScan

Following the outcome of the simulation study which confirmed the relative advantage of NetScan over its conventional counterparts, this section applies the method to real data of street crime incidents to gain insights into the patterns of their concentration along street networks.

### *Study area*

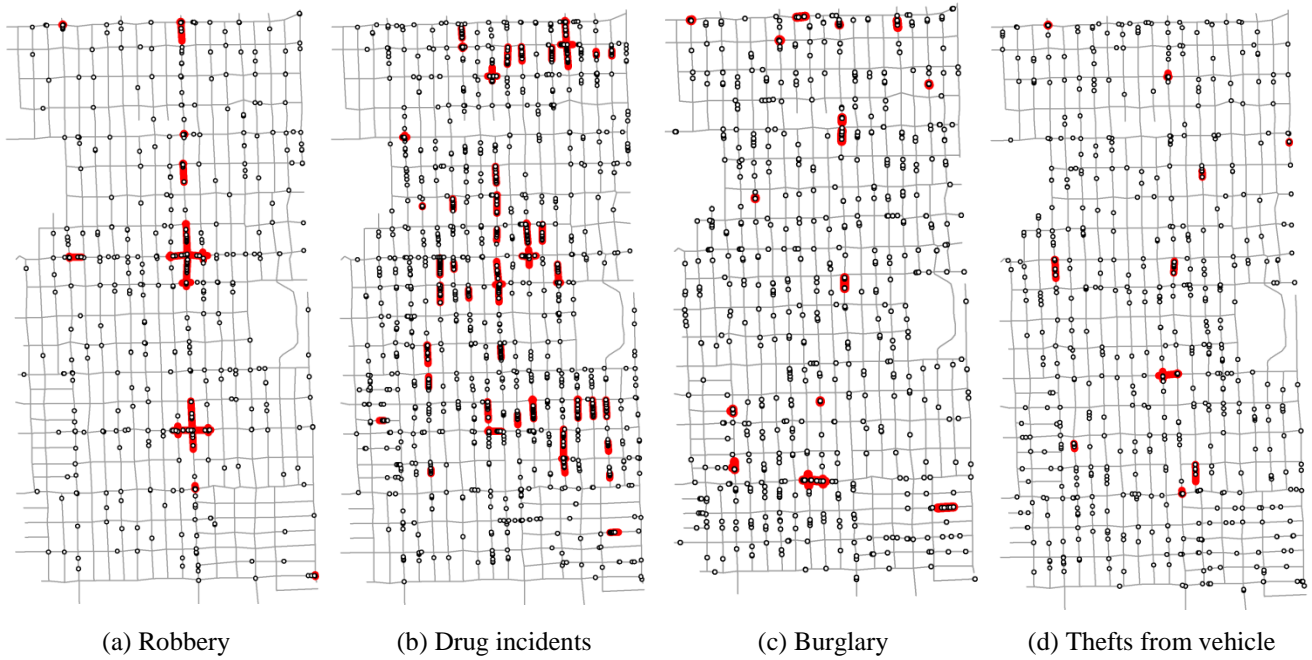
The study area is taken from the Englewood neighbourhood in the mid-southside of Chicago, covering an area of 7,000ft (2,133.6m) by 13,000ft (3,962.4m) with the streets running for a total of 70.13miles (112.86km). The street network is relatively uniform with a major street running at every mile. There is a rapid-transit station in the middle, which forms the centre of the local shopping district.

The study area houses mainly single-family homes and small apartment buildings, many of which are now abandoned. It is known as an area of high risks, with over 40% of its population lives in poverty. The data used are the most common types of street crime: namely, robberies, drug-related incidents, burglaries and thefts from vehicles, recorded by the Chicago Police Department in the year 2000. While all four types of crime saw a high volume of offences during the study period, they are expected to show marked difference in the patterns of their distribution, which should help illustrate how the proposed methods detect hotspots differently. Within the study area, there were 1563 cases of drug-related incidents, 520 robberies, 784 burglaries, and 520 cases of thefts from vehicles recorded in 2000.

### *Application of NetScan for empirical analysis*

Before NetScan can be utilised for hotspot analysis, the size of its search window needs to be considered carefully in relation to the size of the study area as well as the lengths of street segments and the average distance between them. Given that the length of the shorter and the longer edges of each block measure approximately 200ft and 300ft respectively, searches were carried out using a network-based search window with the maximum length of 300ft to ensure a sufficient coverage of each line segment.

Figures 7(a)~7(d) respectively show the crime hotspots detected with NetScan for the four different crime types. The locations of crime incidents are shown with small white circles and the hotspots with bold (red) line segments. While it is difficult to identify the difference in the distribution of incidents between different crime types, there is marked difference in the location of their hotspots.



**Figure 7.** Hotspots of crime incidents in Englewood, Chicago IL in 2000, detected with NT-Scan. Locations of crime incidents and the network-based hotspots are shown respectively with small white circles and bold (red) line segments from the four types of street crimes: (a) robbery, (b) drug-related incidents, (c) burglary, and (d) thefts from vehicle.

For instance, almost all hotspots of robbery incidents are confined to the street running north-south in the centre of the study area and the surrounding shopping district, whereas hotspots of drug-related incidents are dispersed across the study area, covering more streets than any of the other three types of crimes do. This is not surprising, as Englewood was known as an epi-centre of drug-related incidents in 2000. The hotspots of burglary cases are generally scattered across the study area with some degree of concentration on a select few streets, which arguably reflects its tendency for (near-)repeat victimisation. Finally, thefts from vehicle form only few small hotspots and are sporadic in their distribution, which is indicative of the opportunistic nature of the crime.

Interestingly, Ashland seems to attract many cases of robbery, burglary and thefts from vehicle, and hosts several hotspots; but this is not the case with drug-related incidents. Instead, drug incidents are found to form hotspots along streets close to Ashland. It confirms the discrete nature of the crime, where good access and low visibility are preferred by those engaging in drug-related transactions.

## 7. Discussion

In the simulation study, all three methods, circular PL-Scan, elliptic PL-Scan and NetScan, detected hotspots at similar locations but there were also some notable differences. In general, circular PL-Scan over-represented the most likely clusters and secondary clusters by covering too wide an area, thereby merging several different clusters into a single hotspot and also including areas with no crime. This tendency was particularly prevalent when searches were conducted around reference points arranged in grid formation. Using reference points assigned along the street network helped improve the performance of circular PL-Scan in detecting clusters, which showed the importance of using an appropriate set of reference points. Even then, the circular search windows covered either a wider area that over-represented a cluster, or formed a small and under-represented cluster that covered only part of the actual cluster.

Compared to circular PL-Scan, elliptic PL-Scan generally returned better results, as it detected clusters in a narrower and more focused form to follow the linear structure of those clusters. As it was the case with circular PL-Scan, reference points assigned along the street network gave stronger performance, reducing the excess area covered by the elliptic search windows. Nevertheless, elliptic PL-Scan also showed a tendency of incorrectly merging multiple clusters from adjacent streets into a single cluster, and this tendency was even more pronounced as the concentration of those cluster increased.

The relatively low accuracy of circular and elliptic PL-Scans in detecting the size and location of hotspots owes inherently to the nature of their planar search window. Existing search windows, regardless of their shape, primarily take the form of a two-dimensional object which tends to cover a large area around the actual cluster. This inevitably increases the expected number of incidents within the window and, thereby, reduces the likelihood ratio. In contrast, NetScan detects hotspots along a street network without having to include a large excess area. This is not to say that a network-based search window can always provide a perfect match with a cluster, as it too, suffers the risk of slight overshooting along the edges of its search window, as was reflected in the non-perfect PPV and the specificity scores. Nevertheless, using a one-dimensional search window helps keep the margin of error to a minimum, especially when compared to the performance of the PL-Scan techniques as confirmed by the Mann-Whitney's U test.

The comparative analysis also revealed that NetScan is sensitive to local changes in cluster locations; i.e. it reacted to the emergence of a new hotspot immediately when the hotspot turned into the most significant cluster through the injection of additional incidents. This result opens a scope for monitoring micro-scale changes in the risk of crime over time. Overall, the analysis carried out in this

paper demonstrated the effectiveness of NetScan to detect clusters and changes thereof for events recorded along a street network.

The empirical analysis illustrated even more clearly the ability of NetScan to pinpoint the exact shape and location of hotspots. It allowed us to depict the characteristics of each crime type through micro-scale detection of hotspots which, until now, has been difficult, owing to the lack of suitable methods for micro-scale hotspot analysis. The capacity of NetScan to provide an accurate profile of clusters for each crime type should also help inform the relevant criminological theory and, thereby, allow us to improve our understanding of where and how each type of crime forms concentration.

## **8. Conclusion**

This paper proposed a new type of hotspot detection method to describe the micro-space variation of locations of crime incidents at the street-address level. It expands on a widely used hotspot detection technique of Spatial Scan Statistic, to facilitate analysis in the network space. The paper demonstrated its effectiveness in detecting the exact shape and position of hotspots, despite the intricate structure of as well as the distance along street networks. Compared to the existing range of Scan Statistic methods, NetScan provides a more accurate solution to micro-scale hotspot analysis that is also sensitive to changes in hotspots.

In the context of epidemiology, datasets are often aggregated to areal units, and providing “approximate location” of clusters (Kulldorff 2001) is perfectly sensible. However, in a situation where immediate attention and intervention are required at specific locations, such as the examples of crime concentrations used in this paper, it is often crucial to identify the exact locations of clusters and their extent. Applying the conventional types of Scan Statistic at the micro-scale of street addresses yielded less accurate results, and this is where the proposed method of NetScan can be at its most effective. For instance, if the detected cluster covers multiple street locations, it would be impossible to identify whether one of these streets is suffering from an exceedingly high risk, or several streets share moderately high risks. By accurately detecting the extent and intensity of crime hotspots, NetScan facilitates the local law enforcement agencies to take an informed and effective counter measure and, in the long run, assist in fighting crime by design and improving the environment effectively. NetScan also proved effective in detecting changes in cluster locations, and this can help monitor emergence and concentration of crime over a period of time.

Needless to say, the level of accuracy of the analyses is bound by the uncertainty and errors that exist within the data. Such errors are caused mainly by an incorrect or inaccurate assignment of crime incidents to street addresses. While this paper did not attempt to adjust or compensate for these errors, they can be remedied to an extent by the methodological advantage of the proposed methods. This is because the extent of a search window can absorb some of the impact of a possible error in crime locations, provided that the margin of error is smaller than the size of the search window. Still, it does not address the issue in its entirety. An awareness of such limitations, and a proper treatment of uncertainty and errors in the dataset, where possible, remains critical when we analyse micro-scale disaggregate data. As with other data analysis studies, the source data need to be re-examined for their accuracy in the event where a highly unlikely result is produced. This includes the case where a particularly small and highly intense hotspot, a hotdot, was detected as it may be caused by dumping of data to a certain location or by duplicated records.

There are several future directions that can be followed from here. Firstly, the network-based search window technique introduced in this paper can be applied to other subjects and contexts outside the domain of crime and policing. In principle, the idea of network-based analysis would suit inquiries of any datasets that are collected along street networks, and it would cater for a range of applications within the urban and social context as mentioned earlier. In order to facilitate these studies, the authors' group is currently developing a tool that is scheduled to be released, along with sample data and a user manual, for wider use by subject experts who wish to analyse hotspots that emerge along street networks.

Secondly, the proposed method can be extended to incorporate the temporal dimension; i.e. for detecting space-time hotspots that emerge and change their size and density over time. The sensitivity analysis carried out in this paper suggests that NetScan is sufficiently sensitive for the purpose of detecting changes in clusters, and this would allow for the development of a space-time NetScan.

Finally, applying NetScan to a wider range of crime data would help us better understand the background as well as the patterns of crime concentration. According to the crime opportunity theory (e.g. Weisburd, et al., 2009; Braga & Weisburd, 2010; Clarke, 2012), micro-scale hotspots are formed as a manifestation of the small places where criminal opportunities converge. A close examination of the characteristics of the micro-scale hotspot locations, especially in terms of the environmental and situational features associated with those small places, would help understand how these places yield criminal opportunities, and why crimes occur more frequently at these places. The empirical analysis showed that NetScan can clearly depict the difference in the locations and the extent of hotspots by



their crime type. Examining the crime-inducing factors in relation to the newly detected micro places, should help progress both theory and practices around the domain of the geography of crime. While investigating the association of these determinants with the respective places is beyond the scope of this paper, it marks an exciting future direction.

## Funding Support

This research was funded in part by the British Academy/Leverhulme Trust Small Research Grant scheme.

## References

- Aamodt, G., Samuelsen, S.O., & Skrondal, A. (2006). A simulation study of three methods for detecting disease clusters. *International Journal of Health Geographics*, 5, 15pp.
- Amin, R., Ritter, E.K., & Cossette, L. (2012). A geospatial analysis of shark attack rates for the Coast of California: 1994–2010. *Journal of Environment and Ecology*, 3, 246–255.
- Benjamini, Y., & Hochberg, Y. (1995). Controlling the false discovery rate: A practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society, Series B (Methodological)*, 57, 289–300.
- Besag, J., & Newell, J. (1991). The detection of clusters in rare diseases. *Journal of the Royal Statistical Society A*, 154, 143–155.
- Boots, B.N., & Getis, A. (1988). *Point pattern analysis*. Thousand Oaks, CA: Sage Publications.
- Braga, A.A., Papachristos, A.V., & Hureau, D.M. (2010). The concentration and stability of gun violence at micro places in Boston 1980–2008. *Journal of Quantitative Criminology*, 26, 33–53.
- Braga, A.A., & Weisburd, D.L. (2010). Editors' introduction: Empirical evidence on the relevance of place in criminology. *Journal of Quantitative Criminology*, 26, 1–6.
- Clarke, R.V. (2012). Opportunity makes the thief. Really? And so what? *Crime Science*, 1, 3.
- Cliff, A.D., & Ord, J.K. (1981). *Spatial processes: models & applications*, London: Pion.
- Cressie, N. (1991), *Statistics for spatial data*, New York, NY: John Wiley & Sons.
- Dahly, D. (2011). Obesity clustering in Cebu, Philippines: An application of SatScan and the spatial scan statistic. *Journal of Epidemiology and Community Health*, 65, A71.
- DeChello, L.M., & Sheehan, T.J. (2007). Spatial analysis of colorectal cancer incidence and proportion of late-stage in Massachusetts residents: 1995–1998. *International Journal of Health Geographics*, 6, 20.
- Desjardins, M.R., Whiteman, A., Casas, I., & Delmelle, E. (2018). Space-time clusters and co-occurrence of chikungunya and dengue fever in Colombia from 2015 to 2016. *Acta Tropica*, 185, 77–85.
- di Bella, E., Corsi, M., Leporatti, L., & Persico, L. (2017). The spatial configuration of urban crime environments and statistical modeling. *Environment and Planning B*, 44, 647–667.
- Diggle, P.J. (2003). *Statistical analysis of spatial point patterns*. New York, NY: Oxford University Press.
- Duczmal, L., & Assunção, R. (2004). A simulated annealing strategy for the detection of arbitrary shaped spatial clusters. *Computational Statistics and Data Analysis*, 45, 269–286.
- Duczmal, L., Kulldorff, M., & Huang, L. (2006). Evaluation of spatial scan statistics for irregularly shaped clusters. *Journal of Computational and Graphical Statistics*, 15, 428–442.
- Flahaut, B., Mouchart, M., Sanmartin, E., & Thomas, I. (2003). The local spatial autocorrelation and the kernel method for identifying black zones: a comparative approach. *Accident Analysis and Prevention*, 35, 991–1004.

- Forsberg, L. et al. (2005). Distance-based methods for spatial and spatio-temporal surveillance. In: A.B. Lawson, & K. Kleinman (Eds.). *Spatial & Syndromic Surveillance for Public Health, 2nd edition* (pp. 115–131). Chichester: John Wiley & Sons.
- Grubestic, T.H., Wei, R., & Murray, A. T. (2014). Spatial clustering overview and comparison: Accuracy, sensitivity, and computational expense. *Annals of the Association of American Geographers*, 104(6), 1134–1156.
- Helbich, M. (2011). Beyond postsuburbia? Multifunctional service agglomeration in Vienna's urban fringe. *Journal of Economic and Social Geography*, 103, 39–52.
- Henry, K.A., Niu, X., & Boscoe, F.P. (2009). Geographic disparities in colorectal cancer survival. *International Journal of Health Geographics*, 8, 48.
- Huang, L., Kulldorff, M., & Gregorio, D. (2007). A spatial scan statistic for survival data. *Biometrics*, 63, 109–118.
- Jemal, A. et al. (2002). A geographic analysis of prostate cancer mortality in the United States. *International Journal of Cancer*, 101, 168–174.
- Kulldorff, M. (1997). A spatial scan statistic. *Communications in Statistics: Theory and Methods*, 26, 1481–1496.
- Kulldorff, M. (2001). Prospective time periodic geographical disease surveillance using a scan statistic. *Journal of the Royal Statistical Society, Series A (Statistics in Society)*, 164, 61–72.
- Kulldorff, M. (2009). *SatScan User Guide for Version 8.0*, <http://www.satscan.org>.
- Kulldorff, M., Huang, L., Pickle, L., & Duczmal, L. (2006). An elliptic spatial scan statistic. *Statistics in Medicine*, 25, 3929–3943.
- Kulldorff, M., & Nagarwalla, N. (1995). Spatial disease clusters: Detection and inference, *Statistics in Medicine*, 14, 799–810.
- Kulldorff, M., Tango, T., & Park, P.J. (2003). Power comparisons for disease clustering tests. *Computational Statistics and Data Analysis*, 42, 665–684.
- Langen, T.A. et al. (2007). Methodologies for surveying herpetofauna mortality on rural highways. *Journal of Wildlife Management*, 71, 1361–1368.
- Li, M. et al. (2019). Sensitivity of disease cluster detection to spatial scales: an analysis with the spatial scan statistic method. *International Journal of Geographical Information Science*, 33(11), 2125–2152.
- Luquero, F.J. et al. (2011). Cholera epidemic in Guinea-Bissau (2008): The importance of “place”. *PLoS One*, 6, e19005.
- Maheu-Giroux, M., & de Blois, S. (2007). Landscape ecology of *Phragmites Australis* invasion in networks of linear wetlands. *Ecology*, 22, 285–301.
- Nakaya, T., & Yano, K. (2010). Visualising crime clusters in a space–time cube: an exploratory data analysis approach using space–time kernel density estimation and scan statistics. *Transactions in GIS*, 14, 223–239.
- Neill, D.B. (2009). Expectation-based Scan Statistics for monitoring spatial time series data. *International Journal of Forecasting*, 25, 498–517.
- Nordin, J.D. et al. (2005). Simulated anthrax attacks and syndromic surveillance. *Emerging Infectious Diseases*, 11, 1394–1398.
- Okabe, A., & Okunuki, K. (2001). A computational method for estimating the demand of retail stores on a street network and its implementation in GIS. *Transactions in GIS*, 5, 209–220.
- Okabe, A., Okunuki, K., & Shiode, S. (2006). The SANET Toolbox: New methods for network spatial analysis. *Transactions in GIS*, 10, 535–550.
- Okabe, A., & Sugihara, K. (2012). *Spatial analysis along networks: statistical and computational methods*. Chichester: John Wiley & Sons.
- Openshaw, S., Charlton, M., Wymer, C., & Craft, A. (1987). A Mark 1 Geographical Analysis Machine for the automated analysis of point data sets. *International Journal of Geographical Information Systems*, 1, 335–358.
- Osei, F.B., & Duker, A.A. (2008). Spatial dependency of v. cholera prevalence on open space refuse dumps in Kumasi, Ghana: A spatial statistical modeling. *International Journal of Health Geographics*, 7, 62.

- Ozdenerol, E., Williams, B.L., Kang, S.Y., & Magsumbol, M.S. (2005). Comparison of spatial scan statistic and spatial filtering in estimating low birth weight clusters. *International Journal of Health Geographics*, 4, 19.
- Patil, G.P., & Taillie, C. (2004). Upper level set scan statistic for detecting arbitrarily shaped hotspots. *Environmental and Ecological Statistics*, 11, 183–197.
- Ramp, D. et al. (2005). Modelling of wildlife fatality hotspots along the Snowy Mountain Highway in New South Wales, Australia. *Biological Conservation*, 126, 474–490.
- Schmicker, R.H. (2013). An application of SaTScan to evaluate the spatial distribution of corner kick goals in major league soccer. *International Journal of Computer Science in Sport*, 12, 70–79.
- Shiode, S. (2008). Analysis of a distribution of point events using the network-based quadrat method. *Geographical Analysis*, 40, 401–422.
- Shiode, S. (2011). Street-Level spatial scan statistic and STAC for analyzing street crime concentrations. *Transactions in GIS*, 15, 365–383.
- Shiode, S., & Shiode, N. (2009). Detection of hierarchical point agglomerations by the network-based variable clumping method. *International Journal of Geographical Information Science*, 23, 75–92.
- Shiode, S., & Shiode, N. (2013). Network-based space–time search window technique for hotspot detection of street-level crime incidents. *International Journal of Geographical Information Science*, 27, 866–882.
- Shiode, S., Shiode, N., Block, R., & Block, C. (2015). Space–time characteristics of micro-scale crime occurrences: An application of a network-based space–time search window technique for crime incidents in Chicago. *International Journal of Geographical Information Science*, 29, 697–719.
- Song, C., & Kulldorff, M. (2003). Power evaluation of disease clustering tests. *International Journal of Health Geographics*, 2, 1–8.
- Spooner, P., Lunt, I. D., Okabe, A., & Shiode, S. (2004). Spatial analysis of roadside acacia populations on a road network using the network *K*-function. *Landscape Ecology*, 19, 491–499.
- Takahashi, K., Kulldorff, M., Tango, T., & Yih, K. (2008) A flexibly shaped space-time scan statistic for disease outbreak detection and monitoring. *International Journal of Health Geographics*, 7, 14.
- Tango, T., & Takahashi, K. (2005). A flexibly shaped spatial scan statistic for detecting clusters. *International Journal of Health Geographics*, 4, 11pp.
- Tuia, D., et al. (2008). Scan statistics analysis of forest fire clusters. *Communications in Nonlinear Sciences and Numerical Simulations*, 13, 1689–1694.
- Turnbull, B., et al. (1990). Monitoring for clusters of disease: Application to leukemia incidence in Upstate New York. *American Journal of Epidemiology*, 132, 136–143.
- Upton, G.J.G., & Fingleton, B. (1985). *Spatial data analysis by example*. Chichester: John Wiley & Sons.
- Weisburd, D., Bruinsma, G.J.N., & Bernasco, W. (2009). Units of analysis in geographic criminology: Historical development, critical issues, and open questions. In: D. Weisburd, G.J.N. Bruinsma, & W. Bernasco (Eds.). *Putting Crime in Its Place: Units of Analysis in Geographic Criminology* (pp. 3–34). New York, NY: Springer Verlag.
- Weisburd, D., Groff, E.R., & Yang, S-M. (2012). *The criminology of place: street segments and our understanding of the crime problem*. Oxford: Oxford University Press.
- Yamada, I., & Thill, J-C. (2004). Comparison of planar and network *K*-functions in traffic accident analysis. *Journal of Transport Geography*, 12, 149–158.